

# Conservative Contextual Linear Bandits

Abbas Kazerouni<sup>1</sup>, Mohammad Ghavamzadeh<sup>2</sup>, and Benjamin Van Roy<sup>3</sup>

<sup>1</sup>Stanford University, abbask@stanford.edu

<sup>2</sup>Adobe Research, ghavamza@adobe.com

<sup>3</sup>Stanford University, bvr@stanford.edu

November 23, 2016

## 1 Abstract

Safety is a desirable property that can immensely increase the applicability of learning algorithms in real-world decision-making problems. It is much easier for a company to deploy an algorithm that is safe, i.e., guaranteed to perform at least as well as a baseline. In this paper, we study the issue of safety in contextual linear bandits that have application in many different fields including personalized ad recommendation in online marketing. We formulate a notion of safety for this class of algorithms. We develop a safe contextual linear bandit algorithm, called *conservative linear UCB* (CLUCB), that simultaneously minimizes its regret and satisfies the safety constraint, i.e., maintains its performance above a fixed percentage of the performance of a baseline strategy, uniformly over time. We prove an upper-bound on the regret of CLUCB and show that it can be decomposed into two terms: **1)** an upper-bound for the regret of the standard linear UCB algorithm that grows with the time horizon and **2)** a constant (does not grow with the time horizon) term that accounts for the loss of being conservative in order to satisfy the safety constraint. We empirically show that our algorithm is safe and validate our theoretical analysis.

## 2 Introduction

Many problems in science and engineering can be formulated as a decision-making problem under uncertainty. Although many learning algorithms have been developed to find a good policy/strategy for these problems, most of them do not provide any guarantee that their resulting policy will perform well, when it is deployed. This is a major obstacle in using learning algorithms in many different fields, such as online marketing, health sciences, and finance. Therefore, developing learning algorithms with *safety* guarantees can immensely increase the applicability of learning in solving decision problems. A policy generated by a

learning algorithm is considered to be safe, if it is guaranteed to perform at least as well as a baseline. The baseline can be either a baseline value or the performance of a baseline strategy. It is important to note that since the policy is learned from data, and data is often random, the generated policy is a random variable, and thus, the safety guarantees are in high probability.

Safety can be studied in both *offline* and *online* scenarios. In the *offline* case, the algorithm learns the policy from a batch of data, usually generated by the current strategy or recent strategies of the company, and the question is whether the learned policy will perform as well as the current strategy or no worse than a baseline value, when it is deployed. This scenario has been recently studied heavily in both *model-based* (e.g., [7]) and *model-free* (e.g., [3, 13, 14, 12, 11, 10, 6]) settings. In the model-based approach, we first use the batch of data and build a simulator that mimics the behavior of the dynamical system under study (online advertisement, hospital’s ER, financial market), and then use this simulator to generate data and learn the policy. The main challenge here is to have guarantees on the performance of the learned policy, given the error in the simulator. This line of research is closely related to the area of robust learning and control. In the model-free approach, we learn the policy directly from the batch of data, without building a simulator. This line of research is related to off-policy evaluation and control. While the model-free approach is more suitable for problems in which we have access to a large batch of data, such as in online marketing, the model-based approach works better in problems in which data is harder to collect, but instead, we have good knowledge about the underlying dynamical system that allows us to build an accurate simulator.

In the *online* scenario, the algorithm learns a policy while interacting with the real system. Although (reasonable) online algorithms will eventually learn a good or an optimal policy, there is no guarantee for their performance along the way (the performance of their intermediate policies), especially at the very beginning, when they perform a large amount of *exploration*. Thus, in order to guarantee safety in online algorithms, it is important to control their exploration and make it more *conservative*. Consider a manager that allows our learning algorithm runs together with her company’s current strategy (baseline policy), as long as it is safe, i.e., the loss incurred by letting a portion of the traffic handled by our algorithm (instead of by the baseline policy) does not exceed a certain threshold. Although we are confident that our algorithm will eventually perform at least as well as the baseline strategy, it should be able to remain alive (not terminated by the manager) long enough for this to happen. Therefore, we should make it more conservative (less exploratory) in a way not to violate the manager’s safety constraint. This setting has been studied in the multi-armed bandit (MAB) [15]. Wu et al. [15] considered the baseline policy as a fixed arm in MAB, formulated safety using a constraint defined based on the performance of the baseline policy (mean of the baseline arm), and modified the UCB algorithm [2] to satisfy this constraint.

In this paper, we study the notion of safety in *contextual linear bandits*, a setting that has application in many different fields including *online personalized ad recommendation*. We first formulate safety in this setting, as a constraint that must hold *uniformly in time*. Our

goal is to design learning algorithms that minimize regret under the constraint that at any given time, their expected sum of rewards should be above a fixed percentage of the expected sum of rewards of the baseline policy. This fixed percentage depends on the amount of risk that the manager is willing to take. We then propose an algorithm, called *conservative linear UCB* (CLUCB), that satisfies the safety constraint. At each round, CLUCB plays the action suggested by the standard linear UCB (LUCB) algorithm (e.g., [5, 8, 1, 4, 9]), only if it satisfies the safety constraint for the worst choice of the parameter in the confidence set, and plays the action suggested by the baseline policy, otherwise. We also prove an upper-bound for the regret of CLUCB, which can be decomposed into two terms. The first term is an upper-bound on the regret of LUCB that grows at the rate  $\log T \sqrt{T}$ . The second term is constant (does not grow with the horizon  $T$ ) and accounts for the loss of being conservative in order to satisfy the safety constraint. This improves over the regret bound derived in [15] for the MAB setting, where the regret of being conservative grows with time. However, the cost of this improvement for us is to have larger multiplicative constants. Finally, we report experimental results that show CLUCB behaves as expected in practice and validate our theoretical analysis.

### 3 Problem Formulation

In this section, we first review the standard linear bandit setting and then introduce the conservative linear bandit formulation considered in this paper.

#### 3.1 Linear Bandit

In the linear bandit setting, at any time  $t$ , the agent is given a set of (possibly) infinitely many actions/options  $\mathcal{A}_t$ , where each action  $a \in \mathcal{A}_t$  is associated with a feature vector  $\phi_a^t \in \mathbb{R}^d$ . At each round  $t$ , the agent should select an action  $a_t \in \mathcal{A}_t$ . Upon selecting  $a_t$ , the agent observes a random reward  $Y_t$  generated as

$$Y_t = \langle \theta^*, \phi_{a_t}^t \rangle + \eta_t, \quad (1)$$

where  $\theta^* \in \mathbb{R}^d$  is an unknown parameter,  $\langle \theta^*, \phi_{a_t}^t \rangle = r_{a_t}^t$  is the expected reward of action  $a_t$  at time  $t$ , i.e.,  $r_{a_t}^t = \mathbb{E}[Y_t]$ , and  $\eta_t$  is a random noise such that

**Assumption 1.** *Each element  $\eta_t$  of the noise sequence  $\{\eta_t\}_{t=1}^\infty$  is conditionally  $\sigma^2$ -sub-Gaussian, i.e.,*

$$\forall \lambda \in \mathbb{R}, \quad \mathbb{E} \left[ e^{\lambda \eta_t} \mid a_{1:t}, \eta_{1:t-1} \right] \leq \exp \left( \frac{\lambda^2 \sigma^4}{2} \right).$$

The sub-Gaussian assumption automatically implies that  $\mathbb{E}[\eta_t \mid a_{1:t}, \eta_{1:t-1}] = 0$  and  $\mathbf{Var}[\eta_t \mid a_{1:t}, \eta_{1:t-1}] \leq \sigma^4$ .

Note that the above formulation contains time-varying actions set and time-dependent feature vectors for each action, and thus, includes the *linear contextual bandit* setting. In

linear contextual bandit, if we denote by  $x_t$ , the state of the system at time  $t$ , the time-dependent feature vector  $\phi_a^t$  for action  $a$  will be equal to  $\phi(x_t, a)$ , the feature vector of state-action pair  $(x_t, a)$ .

We also make the following standard assumption on  $\theta^*$  and feature vectors:

**Assumption 2.** *There exist  $B, D \geq 0$  such that  $\|\theta^*\|_2 \leq B$  and  $\langle \theta^*, \phi_a^t \rangle \in [0, D]$ , for all  $t$  and all  $a \in \mathcal{A}_t$ .*

We define  $\mathcal{B} = \{\theta \in \mathbb{R}^d : \|\theta\|_2 \leq B\}$  and  $\Phi = \{\phi \in \mathbb{R}^d : \langle \theta^*, \phi \rangle \in [0, D]\}$  to be the parameter space and feature space, respectively.

Obviously, if the learner knows the value of  $\theta^*$ , at each round  $t$ , she will choose the optimal action  $a_t^* = \arg \max_{a \in \mathcal{A}_t} \langle \theta^*, \phi_a^t \rangle$ . Since  $\theta^*$  is unknown, the goal of the learner is to maximize her cumulative expected reward after  $T$  rounds,  $\sum_{t=1}^T \langle \theta^*, \phi_{a_t}^t \rangle$ , or equivalently, to minimize its (pseudo)-regret

$$R_T = \sum_{t=1}^T \langle \theta^*, \phi_{a_t^*}^t \rangle - \sum_{t=1}^T \langle \theta^*, \phi_{a_t}^t \rangle, \quad (2)$$

which is the difference between the sum of expected rewards of the optimal and learner's strategies.

### 3.2 Conservative Linear Bandit

The conservative linear bandit setting is exactly the same as linear bandit, except that there exists a baseline policy  $\pi_b$  (the company's strategy) that at each time  $t$ , selects the action  $b_t \in \mathcal{A}_t$  and incurs the expected reward  $r_{b_t}^t = \langle \theta^*, \phi_{b_t}^t \rangle$ . We assume that the expected reward of the actions taken by the baseline policy,  $r_{b_t}^t$ , are known. This is often a reasonable assumption, since we usually have access to a large amount of data generated by the baseline policy, as it's our company's strategy, and thus, have a good estimate of its performance.

Another difference between the conservative and standard linear bandit settings is the *performance constraint*, which is defined as follows:

**Definition 1** (Performance Constraint). *At each time  $t$ , the difference between the performances of the learner and the baseline policy should not exceed a pre-defined fraction  $\alpha \in (0, 1)$  of the baseline performance. This constraint may be written more formally as*

$$\sum_{i=1}^t r_{b_i}^i - \sum_{i=1}^t r_{a_i}^i \leq \alpha \sum_{i=1}^t r_{b_i}^i,$$

or equivalently as

$$\sum_{i=1}^t r_{a_i}^i \geq (1 - \alpha) \sum_{i=1}^t r_{b_i}^i. \quad (3)$$

The parameter  $\alpha \in (0, 1)$  controls how conservative the learner should be. Small values of  $\alpha$  show that only small losses are tolerated, and thus, the learner should be overly conservative, whereas large values of  $\alpha$  indicate that the manager is willing to take risk, and thus, the learner can explore more and be less conservative. Here given the value of  $\alpha$ , the goal of the learner is to select her action in a way to both minimize her regret (2) and satisfy the performance constraint (3). In the next section, we propose a linear bandit algorithm, called *conservative linear upper confidence bound* (CLUCB), to achieve this goal.

## 4 A Conservative Linear Bandit Algorithm

In this section, we propose a linear bandit algorithm, called CLUCB, that is based on the *optimism in the face of uncertainty* principle, and given the value of  $\alpha$ , both minimizes the regret (2) and satisfies the performance constraint (3). Algorithm 1 contains the pseudocode of CLUCB. At each round  $t$ , CLUCB uses the previous observations and builds a confidence set  $\mathcal{C}_t$  that with high probability contains the unknown parameter  $\theta^*$ . It then selects the *optimistic action*

$$a'_t = \operatorname{argmax}_{a \in \mathcal{A}_t} \max_{\theta \in \mathcal{C}_t} \langle \theta, \phi_a^t \rangle,$$

which has the best performance among all the actions available in  $\mathcal{A}_t$ , within the confidence set  $\mathcal{C}_t$ . In order to make sure that constraint (3) is satisfied, the algorithm plays the optimistic action  $a'_t$ , only if it satisfies the constraint for the worst choice of the parameter  $\theta \in \mathcal{C}_t$ . To make this more precise, let  $S_{t-1}$  be the set of rounds  $i$  before round  $t$  at which CLUCB has played the optimistic action, i.e.,  $a_i = a'_i$ . In other words,  $S_{t-1}^c = \{1, 2, \dots, t-1\} - S_{t-1}$  is the set of rounds  $j$  before round  $t$  at which CLUCB has followed the baseline policy, i.e.,  $a_j = b_j$ .

In order to guarantee that it does not violate constraint (3), at each round  $t$ , CLUCB plays the optimistic action, i.e.,  $a_t = a'_t$ , only if

$$\min_{\theta \in \mathcal{C}_t} \left[ \sum_{i \in S_{t-1}^c} r_{b_i}^i + \left\langle \theta, \overbrace{\sum_{i \in S_{t-1}} \phi_{a_i}^i}^{z_{t-1}} \right\rangle + \langle \theta, \phi_{a'_t}^t \rangle \right] \geq (1 - \alpha) \sum_{i=1}^t r_{b_i}^i,$$

and plays the baseline action, i.e.,  $a_t = b_t$ , otherwise. In the next section, we will describe how CLUCB constructs and updates the confidence sets  $\mathcal{C}_t$ .

### 4.1 Construction of Confidence Sets

Since the expected reward of playing the action suggested by the baseline policy  $\pi_b$  at each round  $t$ , i.e.,  $r_{b_t}^t$ , is known ahead of time, playing a baseline action does not provide any new information about the unknown parameter  $\theta^*$ . Thus, CLUCB initializes its confidence set to  $\mathcal{B}$  and keeps it this way while it plays the baseline action. After a number of rounds when enough budget is saved, CLUCB starts exploring actions other than those suggested by  $\pi_b$  and constructs and updates its confidence set. We follow the approach of [9] to build

---

**Algorithm 1** CLUCB
 

---

**Input:**  $\alpha, \mathcal{A}, \mathcal{B}$   
**Initialize:**  $S_0 = \emptyset$ ,  $z_0 = \mathbf{0} \in \mathbb{R}^d$ , and  $\mathcal{C}_1 = \mathcal{B}$   
**for**  $t = 1, 2, 3, \dots$  **do**  
   Find  $(a'_t, \tilde{\theta}_t) = \arg \max_{(a, \theta) \in \mathcal{A}_t \times \mathcal{C}_t} \langle \theta, \phi_a^t \rangle$   
   Find  $L_t = \min_{\theta \in \mathcal{C}_t} \langle \theta, z_{t-1} + \phi_{a'_t}^t \rangle$   
   **if**  $L_t + \sum_{i \in S_{t-1}^c} r_{b_i}^i \geq (1 - \alpha) \sum_{i=1}^t r_{b_i}^i$  **then**  
     Play  $a_t = a'_t$  and observe reward  $Y_t$  defined by (1)  
     Set  $z_t = z_{t-1} + \phi_{a_t}^t$ ,  $S_t = S_{t-1} \cup t$ , and  $S_t^c = S_{t-1}^c$   
     Given  $(a_t, Y_t)$ , update the confidence set  $\mathcal{C}_{t+1}$  according to (5)  
   **else**  
     Play  $a_t = b_t$  and observe reward  $Y_t$  defined by (1)  
     Set  $z_t = z_{t-1}$ ,  $S_t = S_{t-1}$ ,  $S_t^c = S_{t-1}^c \cup t$ , and  $\mathcal{C}_{t+1} = \mathcal{C}_t$   
   **end if**  
**end for**

---

confidence sets for the unknown parameter  $\theta^*$ . At each round  $t$  that CLUCB plays the optimistic action, we first calculate the least square estimate of the unknown parameter, given the data that have been observed so far  $\{(\phi_{a_i}^i, Y_i)\}_{i \in S_{t-1}}$ , as

$$\hat{\theta}_t = \operatorname{argmin}_{\theta \in \mathcal{B}} \sum_{i \in S_{t-1}} \left( Y_i - \langle \theta, \phi_{a_i}^i \rangle \right)^2, \quad (4)$$

and then update the confidence set as

$$\mathcal{C}_t = \left\{ \theta \in \mathcal{C}_{t-1} : \sqrt{\sum_{i \in S_{t-1}} \left( \langle \theta, \phi_{a_i}^i \rangle - \langle \hat{\theta}_t, \phi_{a_i}^i \rangle \right)^2} \leq \beta(m_{t-1}, \delta) \right\}, \quad (5)$$

where  $m_{t-1} = |S_{t-1}|$  is the number of optimistic actions played prior to round  $t$ ,  $\delta \in (0, 1)$  is the desired confidence level, and

$$\beta(m_t, \delta) = 16d\sigma^2 \log \left( \frac{2(m_{t-1} + 1)}{\delta} \right) + \frac{2}{m_{t-1} + 1} \left( 16D + \sqrt{8\sigma^2 \log \left( \frac{4(m_{t-1} + 1)^2}{\delta} \right)} \right). \quad (6)$$

We also let  $n_{t-1} = |S_{t-1}^c| = t - 1 - m_{t-1}$  denote the number of rounds that the baseline policy has been followed prior to round  $t$ . It is important to note that (5) defines a decreasing sequence of confidence sets, i.e.,  $\mathcal{C}_1 \supseteq \mathcal{C}_2 \supseteq \mathcal{C}_3 \supseteq \dots$ . The following proposition shows that the constructed confidence sets contain the true parameter  $\theta^*$  with high probability.

**Proposition 1.** *For any  $\delta > 0$  and with  $\mathcal{C}_t$  defined in (5), we have*

$$\mathbb{P} \left[ \theta^* \in \mathcal{C}_t, \forall t \in \mathbb{N} \right] \geq 1 - 2\delta. \quad (7)$$

Proposition 1 is a special case of Proposition 6 in [9] for the family of linear functions and we omit its proof here. This proposition indicates that the CLUCB algorithm satisfies the performance constraint (3) at any time  $t$  with probability at least  $1 - 2\delta$ . This is because at any time  $t$ , CLUCB ensures that the constraint (3) holds all  $\theta \in \mathcal{C}_t$ .

## 4.2 Regret Analysis of CLUCB

In the following, we prove a regret bound for the CLUCB algorithm.

Let  $\Delta_{b_t}^t = r_{a_t^*}^t - r_{b_t}^t$  be the difference between the expected rewards of the optimal and baseline actions at time  $t$ . We call  $\Delta_{b_t}^t$  the *baseline gap* at time  $t$ . It indicates how sub-optimal the action suggested by the baseline policy is at time  $t$ . We make the following assumption on the performance of the baseline policy  $\pi_b$ .

**Assumption 3.** *There exist  $0 \leq \Delta_l \leq \Delta_h$  and  $0 < r_l < r_h$  such that at each round  $t$ ,*

$$\Delta_l \leq \Delta_{b_t}^t \leq \Delta_h \quad \text{and} \quad r_l \leq r_{b_t}^t \leq r_h. \quad (8)$$

An obvious candidate for both  $\Delta_h$  and  $r_h$  is  $D$ , as all the mean rewards are confined in  $[0, D]$ . The reward lower-bound  $r_l$  ensures that at each round, the baseline policy maintains a minimum level of performance. Finally,  $\Delta_l = 0$  is a reasonable candidate for the lower-bound of the baseline gap.

The following proposition shows that the regret of CLUCB can be decomposed into the regret of linear UCB (LUCB) and a regret caused by being conservative in order to satisfy the constraint (3).

**Proposition 2.** *The regret of CLUCB can be decomposed into two terms as follows:*

$$R_T(\text{CLUCB}) \leq R_{S_T}(\text{LUCB}) + n_T \Delta_h, \quad (9)$$

where  $R_{S_T}(\text{LUCB})$  is the (pseudo) regret of LUCB at rounds  $t \in S_T$ .

*Proof.* From the definition of regret (2), we have

$$\begin{aligned} R_T(\text{CLUCB}) &= \sum_{t=1}^T r_{a_t^*}^t - \sum_{t=1}^T r_{a_t}^t = \sum_{t \in S_T} r_{a_t^*}^t + \sum_{t \in S_T^c} r_{a_t^*}^t - \left( \sum_{t \in S_T} r_{a_t}^t + \sum_{t \in S_T^c} r_{b_t}^t \right) \\ &= \sum_{t \in S_T} (r_{a_t^*}^t - r_{a_t}^t) + \sum_{t \in S_T^c} (r_{a_t^*}^t - r_{b_t}^t) = \sum_{t \in S_T} (r_{a_t^*}^t - r_{a_t}^t) + \sum_{t \in S_T^c} \Delta_{b_t}^t \\ &\leq \sum_{t \in S_T} (r_{a_t^*}^t - r_{a_t}^t) + n_T \Delta_h. \end{aligned} \quad (10)$$

The result follows by the fact that for  $t \in S_T$ , CLUCB selects the exact same action suggested by LUCB, and thus, the first term in (10) represents LUCB's regret for rounds  $t \in S_T$ .  $\square$

The regret bound of LUCB for the confidence set defined by (5) can be derived from the results in [9]. Therefore, in order to bound the regret of CLUCB, we only need to find an upper bound on  $n_T$ , i.e., the number of times it deviates from LUCB and selects the action suggested by the baseline policy. We start this part of the proof with the following definition.

**Definition 2.** For each  $\phi \in \Phi$  and each  $t \in \{1, \dots, T\}$ , we define the width of the confidence set  $\mathcal{C}_t$  for the feature vector  $\phi$  as

$$w_t(\phi) = \sup_{\theta \in \mathcal{C}_t} \langle \theta, \phi \rangle - \inf_{\theta \in \mathcal{C}_t} \langle \theta, \phi \rangle. \quad (11)$$

The following lemma bounds the sum of the widths of the confidence sets built by CLUCB for the feature vectors of the optimistic actions (those not suggested by the baseline policy) it selects. This lemma plays an essential role in bounding both terms in the regret decomposition (9).

**Lemma 3.** For each  $t \in \{1, \dots, T\}$  and any sequence of the feature vectors of optimistic actions selected by CLUCB up to round  $t$ , i.e.,  $\{\phi_{a_i}^i\}_{i \in S_t}$ , we have

$$\sum_{i \in S_t} w_i(\phi_{a_i}^i) = O\left(Dd\sigma \log\left(\frac{2Bm_t}{\delta}\right)\sqrt{m_t}\right), \quad (12)$$

where  $m_t = |S_t|$  is the cardinality of the set  $S_t$ .

*Proof.* Note that CLUCB updates its confidence set only when it selects an optimistic action, i.e., at rounds  $i \in S_t$ . According to Lemma 5 in [9], we have

$$\sum_{i \in S_t} w_i(\phi_{a_i}^i) \leq 1 + D \dim_E(\mathcal{F}, m_t^{-1}) + 4\sqrt{\dim_E(\mathcal{F}, m_t^{-1})\beta(m_t, \delta)m_t}, \quad (13)$$

where  $\mathcal{F} = \{f_\theta : \Phi \rightarrow \mathbb{R} \mid f_\theta(\phi) = \langle \theta, \phi \rangle, \theta \in \mathcal{B}\}$  and  $\dim_E(\mathcal{F}, \epsilon)$  is the  $\epsilon$ -eluder dimension of  $\mathcal{F}$ . On the other hand, Proposition 11 in [9] suggests that

$$\dim_E(\mathcal{F}, m_t^{-1}) = O\left(d \log(Bm_t)\right). \quad (14)$$

Moreover, from (6), it is easy to see that for sufficiently large  $m_t$ , we have

$$\beta(m_t, \delta) = O\left(d\sigma^2 \log\left(\frac{2m_t}{\delta}\right)\right). \quad (15)$$

The result follows by plugging (14) and (15) into (13).  $\square$

The following proposition bounds the regret of LUCB in the regret decomposition (9). Although this bound is known in the literature, we rederive it here as a direct consequence of Lemma 3. Let  $\mathcal{E}$  be the event that  $\theta^* \in \mathcal{C}_t$ ,  $\forall t \in \mathbb{N}$ , which according to Proposition 1 holds with probability at least  $1 - 2\delta$ .



**Proposition 4.** *On the event  $\mathcal{E}$ , for any horizon  $T \in \mathbb{N}$ , we have*

$$R_{S_T}(\text{LUCB}) = \sum_{t \in S_T} (r_{a_t^*}^t - r_{a_t}^t) \leq 30Dd\sigma \log\left(\frac{2BT}{\delta}\right) \sqrt{T}. \quad (16)$$

*Proof.* For each round  $t \in S_T$ , we define  $\tilde{\theta}_t = \operatorname{argmax}_{\theta \in \mathcal{C}_t} \max_{a \in \mathcal{A}_t} \langle \theta, \phi_a^t \rangle$ . Since  $\theta^* \in \mathcal{C}_t, \forall t \in \mathbb{N}$  on the event  $\mathcal{E}$ , by the definition of  $a_t$  from Algorithm 1 and the definition of  $\tilde{\theta}_t$ , it follows that  $\langle \theta^*, \phi_{a_t^*}^t \rangle \leq \langle \tilde{\theta}_t, \phi_{a_t}^t \rangle$ , for each  $t \in S_T$ . Therefore, we may write

$$\begin{aligned} \sum_{t \in S_T} (r_{a_t^*}^t - r_{a_t}^t) &= \sum_{t \in S_T} \left( \langle \theta^*, \phi_{a_t^*}^t \rangle - \langle \theta^*, \phi_{a_t}^t \rangle \right) \\ &\leq \sum_{t \in S_T} \left( \langle \tilde{\theta}_t, \phi_{a_t}^t \rangle - \langle \theta^*, \phi_{a_t}^t \rangle \right) \leq \sum_{t \in S_T} \left( \sup_{\theta \in \mathcal{C}_t} \langle \theta, \phi_{a_t}^t \rangle - \inf_{\theta \in \mathcal{C}_t} \langle \theta, \phi_{a_t}^t \rangle \right) \\ &= \sum_{t \in S_T} w_t(\phi_{a_t}^t) \leq 30Dd\sigma \log\left(\frac{2Bm_T}{\delta}\right) \sqrt{m_T}, \end{aligned}$$

where the last inequality is the result of Lemma 3. The result follows from the fact that  $m_T \leq T$ .  $\square$

The following theorem provides a bound on the number of rounds at which CLUCB acts conservatively and follows the baseline policy  $\pi_b$ . In [15], such a bound has been derived for the MAB setting which grows logarithmically with time. However, we show in the following theorem that CLUCB plays conservatively only in a finite number of time steps that does not grow with time. Since the MAB setting can be viewed as a special case of our formulation, the same time-independent bound carries out there as well. It is also worth mentioning that although our bound improves over that of [15] in terms of time, but it will incur larger multiplicative constants.

**Theorem 5.** *On the event  $\mathcal{E}$ , for any horizon  $T \in \mathbb{N}$ , we have*

$$n_T = O\left(\frac{D^2 d^2 \sigma^2}{\alpha r_l (\alpha r_l + \Delta_l)} \left[ \log\left(\frac{\sqrt{B/\delta} D d \sigma}{\alpha r_l + \Delta_l}\right) \right]^2\right). \quad (17)$$

*Proof.* Let  $\tau$  be the last round that CLUCB follows the baseline policy (plays conservatively), i.e.,  $\tau = \max\{1 \leq t \leq T \mid a_t = b_t\}$ . From Algorithm 1, at round  $\tau$ , we have

$$\min_{\theta \in \mathcal{C}_\tau} \left\langle \theta, \phi_{a_\tau}^\tau + \sum_{t \in S_{\tau-1}} \phi_{a_t}^t \right\rangle + \sum_{t \in S_{\tau-1}^c} r_{b_t}^t < (1 - \alpha) \sum_{t=1}^{\tau} r_{b_t}^t. \quad (18)$$

On the other hand, since CLUCB confidence sets, defined by (5), are nested (i.e.,  $\mathcal{C}_1 \supseteq \mathcal{C}_2 \supseteq \mathcal{C}_3 \supseteq \dots$ ), we may write

$$\min_{\theta \in \mathcal{C}_\tau} \left\langle \theta, \phi_{a_\tau}^\tau + \sum_{t \in S_{\tau-1}} \phi_{a_t}^t \right\rangle \geq \min_{\theta \in \mathcal{C}_\tau} \langle \theta, \phi_{a_\tau}^\tau \rangle + \sum_{t \in S_{\tau-1}} \min_{\theta \in \mathcal{C}_t} \langle \theta, \phi_{a_t}^t \rangle. \quad (19)$$

Combining (18) and (19), we obtain

$$\min_{\theta \in \mathcal{C}_\tau} \langle \theta, \phi_{a'_\tau}^\tau \rangle + \sum_{t \in S_{\tau-1}} \min_{\theta \in \mathcal{C}_t} \langle \theta, \phi_{a_t}^t \rangle + \sum_{t \in S_{\tau-1}^c} r_{b_t}^t < (1 - \alpha) \sum_{t=1}^{\tau} r_{b_t}^t. \quad (20)$$

Using the fact that  $\tau = m_{\tau-1} + n_{\tau-1} + 1$ , we may rewrite (20) as

$$\begin{aligned} \alpha \sum_{t=1}^{\tau} r_{b_t}^t &< \left[ r_{b_\tau}^\tau - \min_{\theta \in \mathcal{C}_\tau} \langle \theta, \phi_{a'_\tau}^\tau \rangle \right] + \sum_{t \in S_{\tau-1}} \left[ r_{b_t}^t - \min_{\theta \in \mathcal{C}_t} \langle \theta, \phi_{a_t}^t \rangle \right] \\ &= \left[ r_{b_\tau}^\tau - \langle \theta^*, \phi_{a'_\tau}^\tau \rangle + \langle \theta^*, \phi_{a'_\tau}^\tau \rangle - \min_{\theta \in \mathcal{C}_\tau} \langle \theta, \phi_{a'_\tau}^\tau \rangle \right] \\ &\quad + \sum_{t \in S_{\tau-1}} \left[ r_{b_t}^t - \langle \theta^*, \phi_{a_t}^t \rangle + \langle \theta^*, \phi_{a_t}^t \rangle - \min_{\theta \in \mathcal{C}_t} \langle \theta, \phi_{a_t}^t \rangle \right] \\ &\stackrel{\text{(a)}}{\leq} \left[ -\Delta_{b_\tau}^\tau + \max_{\theta \in \mathcal{C}_\tau} \langle \theta, \phi_{a'_\tau}^\tau \rangle - \min_{\theta \in \mathcal{C}_\tau} \langle \theta, \phi_{a'_\tau}^\tau \rangle \right] \\ &\quad + \sum_{t \in S_{\tau-1}} \left[ -\Delta_{b_t}^t + \max_{\theta \in \mathcal{C}_t} \langle \theta, \phi_{a_t}^t \rangle - \min_{\theta \in \mathcal{C}_t} \langle \theta, \phi_{a_t}^t \rangle \right] \\ &= -(m_{\tau-1} + 1)\Delta_l + w_\tau(\phi_{a'_\tau}^\tau) + \sum_{t \in S_{\tau-1}} w_t(\phi_{a_t}^t), \end{aligned} \quad (21)$$

where **(a)** follows from the fact that on event  $\mathcal{E}$ ,  $\theta^* \in \mathcal{C}_t$  for all  $t \in \{1, \dots, T\}$ , and thus,  $\langle \theta^*, \phi_{a'_t}^t \rangle \leq \max_{\theta \in \mathcal{C}_t} \langle \theta, \phi_{a'_t}^t \rangle$  for all  $t \in \{1, \dots, T\}$ . Since  $a_t = a'_t$  for all  $t \in S_T$ , for these  $t$ 's, we also have  $\langle \theta^*, \phi_{a'_t}^t \rangle \leq \max_{\theta \in \mathcal{C}_t} \langle \theta, \phi_{a_t}^t \rangle$ .

From Lemma 3, it follows that

$$w_\tau(\phi_{a'_\tau}^\tau) + \sum_{t \in S_{\tau-1}} w_t(\phi_{a_t}^t) \leq 30Dd\sigma \log \left( \frac{2B(m_{\tau-1} + 1)}{\delta} \right) \sqrt{m_{\tau-1} + 1}. \quad (22)$$

On the other hand, from Assumption 3, we may write

$$\alpha \tau r_l \leq \alpha \sum_{t=1}^{\tau} r_{b_t}^t. \quad (23)$$

Combining (22) and (23) with (21) and rearranging the terms, we have

$$\alpha n_{\tau-1} r_l < -(m_{\tau-1} + 1)(\Delta_l + \alpha r_l) + 30Dd\sigma \log \left( \frac{2B(m_{\tau-1} + 1)}{\delta} \right) \sqrt{m_{\tau-1} + 1}. \quad (24)$$

Since RHS of (24) has a finite upper-bound as it is positive only for a finite range of  $m_{\tau-1}$ , applying the result of Lemma 7 reported in Appendix A, with  $m = m_{\tau-1} + 1$ , gives us

$$\alpha n_{\tau-1} r_l = O \left( \frac{D^2 d^2 \sigma^2}{\alpha r_l + \Delta_l} \left[ \log \left( \frac{\sqrt{B/\delta} D d \sigma}{\alpha r_l + \Delta_l} \right) \right]^2 \right).$$

The result follows from the fact that  $n_T = n_\tau = n_{\tau-1} + 1$ .  $\square$

We now have all the necessary ingredients to derive a regret bound on the performance of the CLUCB algorithm. We report the regret bound of CLUCB in Theorem 6, whose proof is a direct consequence of the results of Propositions 2 and 4, and Theorem 5.

**Theorem 6.** *With probability at least  $1 - 2\delta$ , the CLUCB algorithm satisfies the performance constraint (3) for all  $t \in \mathbb{N}$ , and has the following regret bound on the regret:*

$$R_T(\text{CLUCB}) = 30Dd\sigma\sqrt{T}\log\left(\frac{2BT}{\delta}\right) + K\frac{\Delta_h}{\alpha r_l(\alpha r_l + \Delta_l)}, \quad (25)$$

where  $K$  is a constant that depends only on the parameters of the problem as

$$K = O\left(\left[Dd\sigma\log\left(\frac{\sqrt{B/\delta}Dd\sigma}{\alpha r_l + \Delta_l}\right)\right]^2\right). \quad (26)$$

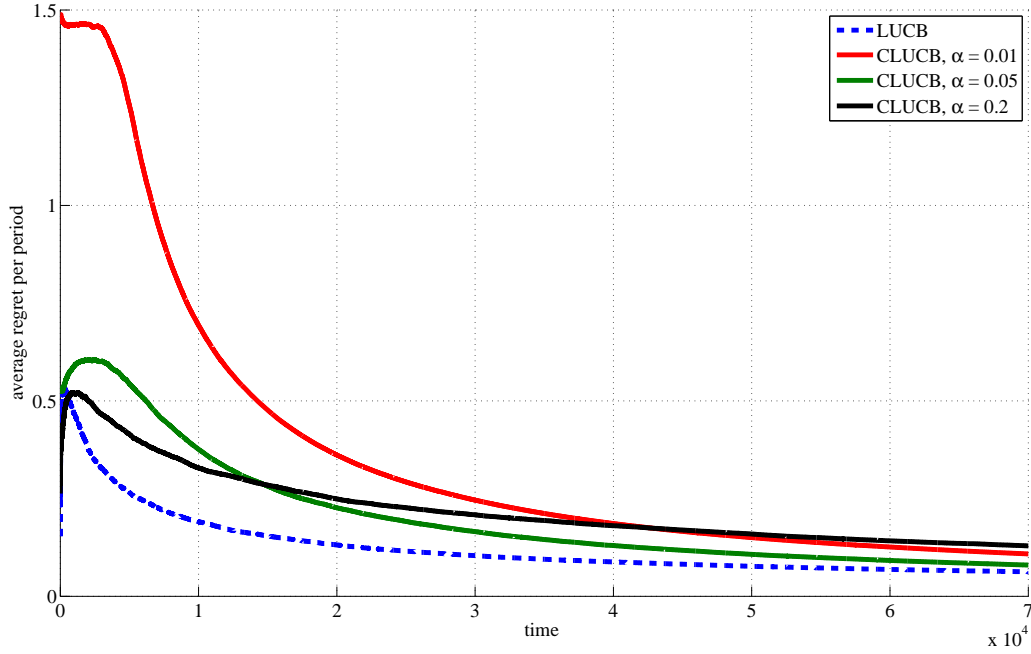
The first term in the regret bound is the regret of LUCB, which grows at the rate of  $\sqrt{T}\log(T)$ . The second term accounts for the loss incurred by being conservative in order to satisfy the performance constraint (3). Our results indicate that this loss does not grow with time (since CLUCB will be conservative only in a finite number of times). This improves over the regret bound derived in [15] for the MAB setting, where the regret of being conservative grows with time. However, as we mentioned earlier, the multiplicative constants in our regret bound are larger than those in the regret bound of [15].

Furthermore, the regret bound of Theorem 6 clearly indicates that CLUCB's regret is larger for smaller values of  $\alpha$ . This perfectly matches the intuition that the agent must be more conservative, and thus, suffers higher regret for smaller values of  $\alpha$ . Theorem 6 also indicates that CLUCB's regret is smaller for smaller values of  $\Delta_h$ , because when the baseline policy  $\pi_b$  is close to optimal, the algorithm does not lose much by being conservative.

## 5 Special Scenarios

In this section, we briefly describe different scenarios which are covered by the problem formulation we made in Section 2.

Our problem formulation in Section 2, allows for time varying action sets as well as time dependent feature vectors for each action. The time varying nature of the action set arises in recommendation and advertisement systems where, as time goes on, new items are available and some old items become obsolete. As an other special case, time dependent feature vectors captures the contextual linear bandit scenario. In contextual bandit, a set of available actions  $\mathcal{A}$  is available each having a feature vector  $\psi_a$ . At any round  $t$ , an observable context  $X_t$  arrives and the expected reward upon playing action  $a$ , is a function of both the context  $X_t$  and the played action  $\psi_a$ . In our model for any action  $a$  and at any round  $t$ , we can build  $\phi_a^t$  as a function of both  $X_t$  and  $\psi_a$  such that the mean reward of playing action  $a$  at time  $t$  given context  $X_t$  is  $\langle \theta^*, \phi_a^t \rangle$ , where  $\theta^*$  is the unknown parameter vector.



**Figure 1:** Average regret of LUCB and CLUCB for different values of  $\alpha$ .

## 6 Simulation Results

In this section, we provide simulation results to illustrate the performance of the proposed CLUCB algorithm. We considered a set of 100 arms each having a feature vector living in  $\mathbb{R}^{10}$  space. Each component of the feature vectors is selected uniformly at random in  $[-1, 1]$  and the parameter  $\theta^*$  is randomly drawn from  $\mathcal{N}(0, 10I_{10})$  such that the mean reward associated to each arm is positive. The noise of the observations is also generated independently from  $\mathcal{N}(0, 4)$  and the mean reward of the base action  $\mu_0$  is set equal to the average of the performances of the second and third best arm.

Fig. 1 depicts the expected reward per period (i.e.,  $\frac{R_t}{t}$ ) of LUCB and CLUCB for different values of  $\alpha$  over a horizon of  $T = 7 \times 10^4$ . Each plot is generated taking average over 20 random realizations of the scenario. As can be seen here, CLUCB plays conservatively at the beginning which incurs a regret larger than LUCB. However, the regret gap between CLUCB and LUCB vanishes over time as CLUCB learns to play the optimal action. As anticipated, Fig. 1 also suggests as  $\alpha$  gets larger the performance of CLUCB converges faster to that of LUCB. On the other hand in our simulations, CLUCB always satisfies the constraints in (3) for all values of  $\alpha$  whereas LUCB violates the safety constraints at an average of 26561 time steps when  $\alpha = 0.01$ . This confirms our theoretical result that CLUCB guarantees the safety constraints at all time while maintaining a regret bound very close to that LUCB.

## References

- [1] Y. Abbasi-Yadkori, D. Pál, and C. Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320, 2011.
- [2] P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning Journal*, 47:235–256, 2002.
- [3] L. Bottou, J. Peters, J. Quinero-Candela, D. Charles, D. Chickering, E. Portugaly, D. Ray, P. Simard, and E. Snelson. Counterfactual reasoning and learning systems: The example of computational advertising. *Journal of Machine Learning Research*, 14:3207–3260, 2013.
- [4] W. Chu, L. Li, L. Reyzin, and R. Schapire. Contextual bandits with linear payoff functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pages 208–214, 2011.
- [5] V. Dani, T. Hayes, and S. Kakade. Stochastic linear optimization under bandit feedback. In *COLT*, pages 355–366, 2008.
- [6] N. Jiang and L. Li. Doubly robust off-policy value evaluation for reinforcement learning. In *Proceedings of the Thirty-Third International Conference on Machine Learning*, 2016.
- [7] M. Petrik, M. Ghavamzadeh, and Y. Chow. Safe policy improvement by minimizing robust baseline regret. In *Advances in Neural Information Processing Systems*, 2016.
- [8] P. Rusmevichientong and J. Tsitsiklis. Linearly parameterized bandits. *Mathematics of Operations Research*, 35(2):395–411, 2010.
- [9] D. Russo and B. Van Roy. Learning to optimize via posterior sampling. *Mathematics of Operations Research*, 39(4):1221–1243, 2014.
- [10] A. Swaminathan and T. Joachims. Batch learning from logged bandit feedback through counterfactual risk minimization. *Journal of Machine Learning Research*, 16:1731–1755, 2015.
- [11] A. Swaminathan and T. Joachims. Counterfactual risk minimization: Learning from logged bandit feedback. In *Proceedings of The 32nd International Conference on Machine Learning*, 2015.
- [12] G. Theodorou, P. Thomas, and M. Ghavamzadeh. Building personal ad recommendation systems for life-time value optimization with guarantees. In *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence*, pages 1806–1812, 2015.
- [13] P. Thomas, G. Theodorou, and M. Ghavamzadeh. High confidence off-policy evaluation. In *Proceedings of the Twenty-Ninth Conference on Artificial Intelligence*, 2015.
- [14] P. Thomas, G. Theodorou, and M. Ghavamzadeh. High confidence policy improvement. In *Proceedings of the Thirty-Second International Conference on Machine Learning*, pages 2380–2388, 2015.
- [15] Y. Wu, R. Shariff, T. Lattimore, and C. Szepesvári. Conservative bandits. In *Proceedings of The 33rd International Conference on Machine Learning*, pages 1254–1262, 2016.

## A Technical Details of the Proof of Theorem 5

In the proof of Theorem 5, we have used the following lemma to bound the RHS of (20).

**Lemma 7.** *For  $m \geq 2$ , the following holds*

$$-m(\Delta_l + \alpha r_l) + 30Dd\sigma \log\left(\frac{2Bm}{\delta}\right)\sqrt{m} = O\left(\frac{D^2d^2\sigma^2}{\alpha r_l + \Delta_l} \left[\log\left(\frac{\sqrt{B/\delta}Dd\sigma}{\alpha r_l + \Delta_l}\right)\right]^2\right). \quad (27)$$

*Proof.* For simplicity in the notation, we let  $c_1 = 30Dd\sigma$ ,  $c_2 = \frac{2B}{\delta}$ , and  $c_3 = (\Delta_l + \alpha r_l)$ , and define the LHS of (27) as a function  $g$  (for  $m \geq 2$ ) of the following form

$$g(m) = -c_3m + c_1\sqrt{m}\log(c_2m).$$

First note that we have

$$g'(m) = -c_3 + \frac{c_1(2 + \log(c_2m))}{2\sqrt{m}}, \quad g''(m) = \frac{-c_1\log(c_2m)}{4m\sqrt{m}},$$

which implies that since  $c_2 > 1$ ,  $g$  is a differentiable concave function over its domain  $[2, \infty)$ . Thus, we can find  $m^*$ , the global maximum of function  $g$ . The first order condition implies that  $g'(m^*) = 0$ , which gives us

$$2 + \log(c_2m^*) = \frac{2c_3}{c_1}\sqrt{m^*}. \quad (28)$$

Plugging this into the definition of  $g$ , gives us

$$g^* = \max_{m \geq 2} g(m) = g(m^*) = c_3m^* - 2c_1\sqrt{m^*}. \quad (29)$$

Now, we take the change of variable  $x = \frac{c_3}{2c_1}\sqrt{m^*}$ , which by (29) gives

$$g^* = \frac{4c_1^2}{c_3}(x^2 - x) \quad (30)$$

On the other hand, (28) becomes

$$2 + \log\left(\frac{4c_2c_1^2}{c_3^2}\right) + 2\log(x) = 4x. \quad (31)$$

Taking exp from both sides gives

$$\frac{e^{4x}}{x^2} = \frac{4c_1^2c_2e^2}{c_3^2}. \quad (32)$$

Now, since  $x^2 \leq e^x$  for all  $x \geq 0$ , it follows from (32) that

$$\frac{4c_1^2c_2e^2}{c_3^2} = \frac{e^{4x}}{x^2} \geq \frac{e^{4x}}{e^x} = e^{3x},$$

which indicates that

$$x \leq \frac{1}{3} \log \left( \frac{4c_1^2 c_2 e^2}{c_3^2} \right).$$

Plugging this into (30) gives that

$$g^* \leq \frac{4c_1^2}{c_3} x^2 \leq \frac{4c_1^2}{9c_3} \left[ \log \left( \frac{4c_1^2 c_2 e^2}{c_3^2} \right) \right]^2 = \frac{4c_1^2}{9c_3} \left[ \log \left( \frac{4c_1^2 c_2 e^2}{c_3^2} \right) \right]^2. \quad (33)$$

Substituting for  $c_1, c_2, c_3$ , we may write

$$g^* = O \left( \frac{D^2 d^2 \sigma^2}{\alpha r_l + \Delta_l} \left[ \log \left( \frac{\sqrt{B/\delta} D d \sigma}{\alpha r_l + \Delta_l} \right) \right]^2 \right). \quad (34)$$

The statement follows from the fact that  $g(m) \leq g^*$  for any  $m \geq 2$ .

□